

This Page Is Inserted by IFW Operations
and is not a part of the Official Record

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

IMAGES ARE BEST AVAILABLE COPY.

**As rescanning documents *will not* correct images,
please do not report the images to the
Image Problem Mailbox.**

1
2
3
4
5
6
7
8
9
0
1
2
3
4
5
6
7
8
9
0
1
2
3
4
5
6
7
8

2
3

4

5

6
7

8

10
11
12
13
14
15

17
18
19
20

21

22
23
24
25
26
27
28

1 optimization (52), abstract interpretation (59,68), instruction specialization (55), and control flow
2 graph reduction (44).

3 The process of producing an optimized version of the original code (30) automatically
4 suppresses some features that can be a hindrance to human malicious code analysis, like
5 overlapping instructions and cast-away branches.
6

7 Optimization (40) is an original way of dealing with polymorphic (10) and other malicious
8 code. The unique ability of optimization (40) to simplify tangled metamorphic code (20) into a
9 readable form can be a crucial advantage in the response to a fast-spreading metamorphic worm
10 (20).

11 **Brief Description of the Drawings**

12 These and other more detailed and specific objects and features of the present invention
13 are more fully disclosed in the following specification, reference being had to the accompanying
14 drawings, in which:
15

16 Figure 1 is an illustration of polymorphic malicious computer code 10.

17 Figure 2 is an illustration of metamorphic malicious computer code 20.

18 Figure 3 is an illustration of apparatus suitable for carrying out the present invention.

19 Figure 4 is an illustration of a method embodiment of the present invention.

20 Figure 5 is an illustration of forward pass steps 42 within the method illustrated in Figure
21 4.
22

23 Figure 6 is an illustration of backward pass steps 43 within the method illustrated in
24 Figure 4.

25 Figure 7 is an example of a Directed Acyclic Graph (DAG).

26 Figure 8 is an example of a control flow graph.

27 Figure 9(a) is a control flow graph for an exemplary section of code before reduction.
28

1 Figure 9(b) is a control flow graph illustrating the code of Figure 9(a) after it has been
2 reduced.

3 Detailed Description of the Preferred Embodiments

4 As used throughout the following specification including claims, the following terms have
5 the following meanings:
6

7 "Malicious computer code" or "malicious code" is any code that is present in a computer
8 without the knowledge and/or without the consent of an authorized user of the computer, and/or
9 any code that can harm the computer or its contents. Thus, malicious code includes viruses,
10 worms, Trojan horses, spam, and adware. At certain places herein, the word "virus" is used
11 generically to include worms and Trojan horses, as well as viruses in the narrow sense.
12

13 "Polymorphic" malicious code is code containing one or more decryption loops and an
14 encrypted virus body that is constant once decrypted

15 "Metamorphic" malicious code is code having a non-constant virus body. Metamorphic
16 code may or may not have decryption loops.

17 "Decryption loop" is a section of malicious code containing instructions to decrypt an
18 encrypted body of the malicious code. The term "decryptor" is often used synonymously with
19 "decryption loop", and sometimes used slightly more generically than "decryption loop".
20

21 "Body" or "virus body" of malicious code is that section of the malicious code that
22 performs the malicious purposes of the code.

23 "Pattern matching" is a technique for recognizing malicious code by looking for patterns
24 or sequences of bits (e.g., signatures) within the code.

25 "Coupled" means any direct or indirect communicative relationship.

26 All of the modules illustrated herein, such as modules 31-36 and 38 illustrated in Figure 3,
27 can be implemented in software, hardware, firmware, and/or any combination thereof. When
28

1 implemented in software, these modules can reside on any computer-readable medium or media
2 such as a hard disk, floppy disk, optical disk, etc.

3 A method embodiment of the present invention determines whether computer code 30
4 contains malicious code. The method comprises the steps of optimizing 40 the computer code 30
5 to produce optimized code; and subjecting the optimized code to a malicious code detection
6 protocol. The malicious code detection protocol can be any protocol for detecting malicious
7 code. Thus, the protocol can be pattern matching, emulation, checksumming, heuristics, tracing,
8 X-raying, algorithmic scanning, or any combination thereof. "Algorithmic scanning" is the use of
9 any custom designed algorithm by which one searches for malicious code. The optimizing 40
10 comprises performing one or more of the following techniques: constant folding 53, copy
11 propagation 54, non-obvious dead code elimination 62,63, peephole optimization 52, code motion
12 49, abstract interpretation 59,68, instruction specialization 55, and control flow graph reduction
13 44. Two or more of these techniques may be combined synergistically.

14 The invention has particular applicability to computer code 30 that is polymorphic 10 or
15 metamorphic 20. When the code 30 is polymorphic 10, in one embodiment the optimizing step
16 40 comprises optimizing just the decryption loop 11, or possibly several decryption loops 11 if
17 the malicious code 10 employs several encryption layers. This is because the viral body 12 is
18 normally written in an already optimal form by the creator of the malicious code 10.

19 When the computer code 30 comprises a decryption loop 11,21 and a viral body 12, 22,
20 one method embodiment of the present invention comprises the steps of optimizing 40 the
21 decryption loop 11,21 to produce optimized loop code; performing a malicious code detection
22 procedure on the optimized loop code; optimizing the body 12, 22 to produce optimized body
23 code; and subjecting the optimized body code to a malicious code detection protocol. This
24 embodiment is particularly useful when the computer code is metamorphic 20. When the
25
26
27
28

1 computer code 30 comprises more than one decryption loop 11, 21, one method embodiment of
2 the present invention comprises the steps of optimizing 40 the outermost decryption loop 11,21 to
3 produce optimized loop code; performing a malicious code detection procedure on the optimized
4 loop code; decrypting the outermost layer, for instance by emulating the optimized loop code;
5 then proceeding in the same way for the second decryption loop, third decryption loop, etc... and
6 all the following innermost encryption layers, until the body 12, 22 is decrypted; optimizing the
7 body 12, 22 to produce optimized body code; and subjecting the optimized body code to a
8 malicious code detection protocol. The malicious code detection procedure can be pattern
9 matching, emulation, checksumming, heuristics, tracing, or algorithmic scanning. The malicious
10 code detection protocol can be pattern matching, emulation, checksumming, heuristics, tracing,
11 X-raying, or algorithmic scanning. The step of optimizing the body can entail using one or more
12 outputs from the step of optimizing the decryption loop and/or the step of performing a malicious
13 code detection procedure on the optimized loop code. When the step of performing a malicious
14 code detection procedure on the optimized loop code indicates that the analyzed code 30 contains
15 malicious code, the steps of optimizing the body and subjecting the optimized body code to a
16 malicious code detection protocol can be aborted. The method can comprise the additional step
17 of revealing encrypted body code. This can be done by emulation or by applying a key gleaned
18 from the optimized loop code.
19
20
21

22 **I. Optimization techniques and their application to polymorphic code 10 and other code 30**
23 **that may contain malicious code.**

24 In this section, we look at specific optimization techniques usable in the present invention,
25 and see how each one of them can be applied to the simplification of polymorphic 10 and other
26 code.
27
28

1 In the following paragraphs, we use two notations for code. One is the classic three-
2 address statement notation often used to describe intermediate code produced in compilers. For
3 instance, the statement:

4 `x := y + z`

5 performs the addition of variables y and z and stores the result in variable x.
6

7 We also use the Intel syntax for x86 microprocessor assembly code. For instance the
8 instruction:

9 `add eax, ebx`

10 performs the addition of registers eax and ebx, stores the result in register eax, and sets the
11 processor flags accordingly. (Note that the left operand is the destination.) When using the term
12 “instruction” within this specification, we refer to processor instructions from the Intel x86
13 instruction set.
14

15 *Uses and definitions*

16 Before proceeding to look into optimization techniques, it is useful to start with the
17 definitions of some common terms.

18 The “uses” of a statement or instruction are the variables whose values are used when the
19 statement or instruction is executed. The “definitions” are the variables whose values are
20 modified when the statement is executed. Variables include registers, processor flags, and
21 memory locations.
22

23 For instance, the statement:

24 `x := y + z`

25 uses variables y and z, and defines variable x. We also say that the statement “kills” any previous
26 definitions of variable x.

27 The x86 instruction:
28

```
1 add    eax, ebx
```

2 uses registers `eax` and `ebx`, and defines register `eax` as well as the overflow, sign, zero, carry,
3 parity, and auxiliary carry flags of the processor. Notice that although the alteration of the flags is
4 just a side effect of the addition, the flags are listed in the definitions set of the instruction.

The instruction:

```
7 | mov byte [edi+esi], 3
```

8 uses registers `edi` and `esi`, and defines whatever memory location the effective address “`edi+esi`”
9 points to. (Note that even though registers `esi` and `edi` appear in the destination operand of the
0 `mov` instruction, they are used and not defined.) Depending on the context, we may be able
1 specify the exact memory location that this instruction defines, or we may have to do a
2 conservative estimate of its definitions set.

Control flow and basic blocks

15 The control flow of a program describes the possible paths it can go along when it is
16 executed. If an execution of a program reaches a conditional branch, such as the “jz” instruction
17 in the following case:

```
18 | label_0:
```

```
19      inc esi
```

```
20      cmp     esi, 10
```

```

21      jz     label 2
22

```

```
23 | label_1:
```

```
24      add    esi, 3
```

```
25 | label 2:
```

```
26      mov     edi, esi
```

27 ret

1 (example 1)

2 This is graphically illustrated in Figure 8. On this control flow graph, the nodes represent
3 instructions or group of instructions; and the degrees represent all possible execution paths.

4 The conditional jump “jz” can be taken or not, depending on the value of register esi. We
5 say that the control flow diverges.

6 We define a basic block as a contiguous set of instructions not interrupted by a branch or
7 the destination of a branch. In the example above, there are three basic blocks: The three
8 instructions between “label_0” and “label_1” form a basic block, so does the single instruction
9 between “label_1” and “label_2”, and so do the two instructions after “label_2.” We often use the
10 term “block” instead of “basic block” in the following text.

11 The successors of a basic block B are the blocks to which control may flow immediately
12 after leaving block B . The predecessors are defined in a similar manner.

13 *Live and dead sets*

14 We say that a variable is live at one point in the program if its value can be used later on
15 during the execution of the program. Otherwise, we say that the variable is dead.

16 For instance, in the example above (example 1), register esi is live on entry into the
17 second basic block, that is at point “label_1,” because its value is used in the execution of the
18 instruction “add esi, 3.” On the other hand, register edi is dead at “label_1,” because its value can
19 never be used before it is defined by the instruction “mov edi, esi.”

20 From the set of live variables at the end of a basic block, it is possible to derive the set of
21 live variables at the beginning of the block by working our way up through the instructions of the
22 block, from the last one to the first one, and applying repeatedly the following data-flow equation.
23 If an instruction I uses the set of variables U and defines the set of variables D , the relation
24 between the live set on entry into I and the live set on exit from I is given by the equation:
25

1 $Live\ set\ on\ entry = (Live\ set\ on\ exit - D) \cup U$

2 In other words, a variable is live before the instruction if it is either used by the
3 instruction, or not killed by the instruction and live after the instruction.

4 Another data-flow equation gives the relation between live variables sets across basic
5 blocks. If block B has successors $S1, S2, \dots, Sn$, then the live set on exit from B is the union of the
6 live sets on entry into the Si 's.

7 $Live\ set\ on\ exit\ from\ block = \cup\ over\ all\ successors\ Si\ (Live\ set\ on\ entry\ into\ Si)$

8 In other words, a variable is live on exit from a block if it is live on entry into at least one
9 successor of the block.

10 Most of the time, the live sets can be computed in linear time, in less than three passes for
11 typical programs.

12 *Dead code elimination*

13 If the definitions set of an instruction contains only dead variables at the point after the
14 instruction, we say that the instruction itself is dead. In such a case, the instruction can be
15 removed from the program without changing the meaning of the program.

16 This transformation is named "dead code elimination". Why would a program contain
17 dead code? Dead code may result from high-level constructs if the programmer overlooked an
18 unneeded variable assignment, but it also very often appears as the result of other optimization
19 techniques we will describe shortly.

20 In polymorphic code 10 produced by viruses, dead code is commonplace. For instance,
21 consider the following snippet of code from a polymorphic decryptor 11 generated by
22 Win32/Junkcomp.

```
23       lea ecx, ds:0ABC5E94Fh  
24  
25       dec cl
```

```

1  sub  al, 0CEh
2  lea  edx, ds:0A979D43Ch
3  inc  cl
4  or   al, 0AFh
5  lea  ebp, ds:0BF8E8B60h
6
7  or   bl, 0B5h
8  bsf  ebx, eax
9  mov  edi, 0B4FA9CF7h
10 rcr  dh, 4Eh
11 bts  edi, ebx
12 imul ebx, esi, 68F2BD76h
13
14 mov  ecx, 0D6FC939Eh

```

Since the last instruction defines register ecx, and ecx is used nowhere in the code before this last definition, the three previous instructions defining ecx or cl are good candidates for dead code elimination. The only catch is that they may also define flags, so we must verify that the flags are also dead after these instructions before we can safely remove them. "lea" does not touch the flags. The flags from "dec cl" are killed by the following "sub" and those from "inc cl" are killed by the following "or". Therefore, it is safe to eliminate these instructions.

The benefits from dead code elimination are numerous. Suppose the instruction stream above is part of a decryption loop 11, and the loop 11 has to be emulated to decrypt the virus body 12. Removing the dead instructions from the loop 11 and then emulating the resulting, simpler code makes the emulation faster. Dead code elimination itself has a cost, but the savings easily outweigh the cost in most cases, since dead code elimination takes place only once, whereas the removed loop instructions might have been executed thousands of times.

1 As used herein, "non-obvious" dead code elimination means removing dead code other
2 than a nop ("no operation") or a simple operation such as cli, sti, clc, stc and others commonly
3 used as single-instruction nop's.

4 Note that emulation of optimized code is slightly different from regular emulation, as the
5 interpreted instructions are not fetched from the emulated memory. Instead, they are fetched from
6 a structure 38 unrelated to the memory that holds symbolic representations of processor
7 instructions, typically a set of nodes in the shape of a control flow graph (see Fig. 3). The
8 optimized instructions may not even have a binary representation. The advantage of this approach
9 is that the memory holding the original code remains unchanged, and the decryption process
10 works even if the bytes of the decryptor 11,21 themselves are used as a decryption key, as is the
11 case in some viruses 10,20.
12

13 If the detection algorithm for the virus is based on loop 11,21 recognition, dead code
14 elimination helps too, by removing unneeded or redundant instructions, thus exposing the more
15 meaningful parts of the code for easier pattern matching. (See the Win32/Dislex example of
16 Illustration E below.) Characteristics of the eliminated instructions, such as the statistic
17 distribution of opcodes in dead code, may also be used for detection.
18

19 Another benefit of dead code elimination is that it may eliminate some anti-emulation
20 code designed to stop antivirus programs. The following snippet of code is taken from the
21 decryptor of Win32/Hezhi.A.
22

```
23 push edx  
24 push edx  
25 lar  edx,  eax  
26 pop  edx  
27 popf  
28
```

1 The “lar” instruction is a rarely used instruction that loads the access rights of a descriptor
2 into a register and modifies the zero flag of the processor. Its presence in the decryptor of the
3 virus is destined to cause some emulators to stop, since they may not know how to emulate the
4 instruction correctly. However, since both edx and the zero flag are dead on exit from the
5 instruction, the “lar” could be discarded as dead code, and the emulation of the optimized code
6 could take place even without proper support for this esoteric instruction.
7

8 Fake import calls may also be eliminated this way if their return values are dead and they
9 have no side effects. (This is unfortunately not the case for Win95/Drill, since it uses the return
10 values of its fake calls to GetModuleHandle, GetTickCount, and other win32 APIs.)

11 *Constant folding*

12 Constant folding consists in replacing expressions that involve only constants by their
13 calculated results, to avoid evaluating them at run time. For instance, the following high-level
14 language statement lends itself to constant folding.
15

16 `i = 1000 + 2 * 3`

17 Rather than generating the code for the multiplication and the addition, a clever compiler
18 will evaluate the value of the expression on the right-hand side of the statement at compile time
19 and generate code for this simple assignment instead:
20

21 `i = 1006`

22 In the context of assembly language, expressions are not apparent, but the idea is the
23 same. Constant folding consists in replacing occurrences of a variable that is known to assume a
24 constant value with the value itself.

25 The following assembly code taken from a sample of Win32/Zmist.A serves to illustrate
26 the transformation:

27 `xor eax, eax`
28

1 sub eax, 87868600

2 push eax

3 After the “xor,” register eax holds the value 0. After the “sub,” eax holds the value
4 78797a00. Thus, we can replace the occurrence of variable eax in the “push” instruction with its
5 constant value at this point, and rewrite the code as:

6
7 xor eax, eax

8 sub eax, 87868600

9 push 78797a00

10 In doing so, we remove register eax from the uses of the “push” instruction, which may
11 have the side effect of exposing dead code. This is an example of the synergy mentioned above.
12 Suppose register eax and the flags defined by the “sub” are dead after the “push.” We could then
13 get rid of the “xor” and the “sub” by dead code elimination.

14
15 The process of constant folding is very similar to emulation. Evaluating an expression
16 written in assembly language is essentially equivalent to performing a partial emulation of the
17 instructions involved in computing the expression.

18 It is a common feature of many polymorphic viruses 10 (and metamorphic viruses 20) to
19 avoid direct use of constants by replacing them with series of instructions producing the desired
20 result. The absence of constants such as looping factors, memory addresses, and decryption keys
21 makes the detection of polymorphic decryptors 11 more difficult. Constant folding can help
22 recover these features.

23
24 To illustrate the benefits of constant folding further, let us use an example related to
25 heuristic detection. Suppose a heuristic engine attempts to detect viral-looking code by searching
26 for small suspicious code snippets. One such snippet may be:

27 cmp word [???+18], 10b
28

1 jnz ???

2 (example 2)

3 This piece of code may appear in the infection routine of viruses that check the COFF
4 signature field at offset 18 (hexadecimal) of the PE header before infecting a file. The question
5 marks designate wildcards for a base register and a branch destination.
6

7 A common anti-heuristic trick for a virus would be to use a slight variant of the code with
8 an equivalent meaning but a different signature such as:

9 mov ax, 10a

10 inc ax ; ax now holds value 10b

11 cmp word [ebx+18], ax

12 jnz dont_infect
13

14 Similar tricks have been played against TBSscan in the past.

15 By applying the constant folding transformation described above and then applying the
16 heuristics to the optimized code, the anti-heuristic trick can be circumvented.

17 *Copy propagation*

18 When a program statement moves the value of a variable into another variable, we say it
19 creates a copy of the variable. The copy is valid as long as both variables remain unchanged.
20

21 For instance, consider the following statements:

22 x := y

23 z := u + x

24 y := u + z

25 x := y + v

26 The first statement creates a copy of variable y into variable x. The third statement
27 invalidates the copy, because variable y is redefined.
28

Copy propagation consists in replacing the variables that are copies of other variables with the originals. In the example above, copy propagation yields the following result:

```
x := y
z := u + y
y := u + z
x := y + v
```

The instance of variable x in the original second statement has been replaced with y, of which it is a copy.

Like constant folding, copy propagation can create new opportunities for dead code elimination. This is another example of the synergy mentioned above. In this example, after removing the reference to variable x in the second statement, the first statement becomes dead code.

In polymorphic code 10, copies are often redundant and can be eliminated. This makes the code 10 clearer to read, easier to parse, and faster to emulate. Look at these few instructions generated by Win32/Simile.A as part of its polymorphic decryptor 11:

```
mov ecx, dword [esi+4000e000]
mov dword [40023ee2], ecx
push dword [40023ee2]
pop dword [40024142]
push dword [40024142]
pop dword [40023c60]
xor dword [40023c60], 8a00e5ca
```

All the first six instructions do is move a value around before it is finally decrypted by the "xor".

1 After copy propagation, the code becomes:

```
2   mov   dword [40023c60], dword [esi+4000e000]
3   xor   dword [40023c60], 8a00e5ca
```

4 This is both easier to understand and faster to emulate. (The double memory-addressing
5 mode of the “mov” is a natural extension of the x86 instruction set.)

6 Notice that copy propagation should not be done for destination operands. The original
7 code is not equivalent to the following instruction!

```
9   xor   dword [esi+4000e000], 8a00e5ca
```

10 *Code motion*

11 One of the goals of optimizing compilers is to produce better code for the parts of a
12 program that are going to be executed the most often. In the absence of programmer hints, it is
13 reasonable enough to attempt optimizing loops the most, especially inner loops.

14 One way to achieve faster loop execution is to move the computation of values that do not
15 change across iterations (so called loop invariants) outside of the loop. For example, assume the
16 following instructions form a decryption loop 11,21:

17 decrypt :

```
19       mov   ebp, [key]
20       xor   [esi], ebp
21       add   esi, 4
22       loop  decrypt
```

23 If we can prove that the memory location holding the key is not affected by the “xor,” we
24 know that register ebp will assume the same value on each loop iteration. Therefore, we can place
25 the initialization of ebp before the loop like this:

```
27       mov   ebp, [key]
```

28

```
1 decrypt :  
2     xor    [esi], ebp  
3     add    esi, 4  
4     loop   decrypt
```

5 The resulting loop has three instructions instead of four, so it will be faster to emulate.

6 Moving computations earlier in the control flow is a common type of code motion, but it
7 is not the only one. Some other similar transformations delay the execution of statements, and
8 possibly duplicate statements, also in an attempt to improve the code in loops.

9 Here we do not discuss the recognition of loops or the exact conditions to use code motion
10 safely. It is enough to rely on the intuitive idea of a loop to see the value of the code motion
11 transformation above.

12 *Peephole optimization*

13 A peephole optimizer 31 is a component that looks at the input stream of machine
14 instructions 30 and makes opportunistic modifications to the stream 30 by removing, replacing, or
15 combining instructions. The peephole optimizer 31 does not know about the meaning of the code
16 30. It just makes simple transformations based on a low-level view of the code 30.

17 The peephole optimizer 31 typically knows a lot about the target architecture, so it can
18 take advantage of special addressing modes and other machine idioms. It may also get rid of
19 back-to-back stores and loads of the same variable, and implement some simple algebraic
20 identities.

21 When dealing with polymorphic code 10, a peephole optimizer 31 can be very useful as
22 the first step 52 of the optimization process 40, as part of an instruction decoder. Polymorphic
23 code 10 is often littered with small sequences of instructions that cancel each other, such as back-
24 to-back negations, complements, or an increment followed by a decrement.

1 Consider a typical example (taken from Win32/Hezhi):

2 rol edx, 1

3 ror edx, 1

4 The two rotations cancel each other. When the peephole optimizer 31 reaches the location
5 of the “rol,” it can look-ahead by one instruction and see that the next instruction is a “ror” of the
6 same register by the same amount, and return a “nop” instead of the “rol.” However, doing this
7 implies an implicit assumption that the flags set by “ror” are dead on exit from the “ror.” This
8 must be carefully verified, either by doing some limited live variable analysis before validating
9 the peephole optimization 52, or by guessing that the flags are dead, and verifying it later in the
10 instruction decoding process. If the assumption about the dead flags turns out to be false, the
11 optimization 52 has to be reversed.

12 Note that this optimization 52 should not preclude the “ror” instruction from being
13 decoded separately at the beginning of a new basic block later on, if it turns out to be the
14 destination of a branch. This peephole optimization 52 is for the instruction sequence starting at
15 the “rol” instruction.

16 A useful peephole optimization 52 is the transformation of push/pop sequences into mov’s
17 (see Win32/Simile example in Illustration F below). This removes the dependency on the stack
18 and introduces more optimization opportunities. However, it can be risky to transform code this
19 way in some contexts, as we will see in detail in a later section.

20 Many similar peephole optimizer 31 tricks can be played, and these will be apparent to
21 people who have some experience working with polymorphic viruses 10. One other case deserves
22 special mention though, the case of back-to-back conditional branches.

23 Two contiguous conditional jumps to the same location that test for complementary
24 conditions (like a jz/jnz pair) can be replaced with one unconditional jump. In a pair of two
25 conditions (like a jz/jnz pair) can be replaced with one unconditional jump. In a pair of two

1 contiguous conditional jumps that test for complementary conditions but have different
2 destinations, the second jump can be replaced with an unconditional jump. Jumps with zero
3 offsets can be replaced with nops. These transformations are all simple, but they are very useful
4 because they simplify the control flow of the code 30.

5
6 In some cases, peephole optimization 52 over a long sequence of instructions might be
7 necessary (for instance for nested push/pop pairs). Implementing the peephole optimizer 31 as a
8 shift-reduce parser helps.

9 *Local vs. global optimization*

10 An optimization is said to be local if it is done at the level of a basic block. It is said to be
11 global if it uses information propagated across basic blocks boundaries. Dead code elimination,
12 constant folding, and copy propagation can all be done locally or globally.

13
14 Local optimizations are less costly and can typically be done in linear time. Most
15 interesting global data-flow problems are proven to be NP-complete, but there is empirical
16 evidence that some can be solved by fast algorithms, at least for programs with a usual control
17 flow structure (and, in this context, polymorphic code 10 does have a usual structure!).

18 In the examples of polymorphic code 10 optimization presented in the Illustrations that are
19 given below, almost all the transformations that were used were local ones, and they gave very
20 good results. Global dead code elimination 63 was the only global optimization implemented, and
21 it brought marginal improvement over local dead code elimination 62.

22
23 It should be noted, however, that two tricks were used to boost local optimizations without
24 paying the extra cost in complexity associated with global optimizations. First, unconditional
25 branches to blocks with only one predecessor were eliminated. This technique is sometimes
26 called “jump removal”, and defeats a common type of polymorphism that consists in slicing the
27
28

code to obfuscate it into little pieces linked together by jumps (see for instance Illustration A on Win32/Zperm.)

Secondly, conditional branches whose conditions fell prey to local optimizations were replaced with jumps or nop's (depending if the branch is always or never taken). Look at this example produced by Win32/Simile.A:

```
mov  dword [4002372a], esi
cmp  esi, dword [4002372a]
jnz  4000b2d9
```

The comparison must always succeed, so the jump is never taken. After copy propagation and instruction specialization, this code became:

```
mov  dword [4002372a], esi
cmp  0, 0
nop
```

Ripe for dead code elimination once the flags of the "cmp" are proven unused.

Abstract interpretation

Abstract interpretation, also called abstract debugging, can be a powerful technique. It consists in modeling the behavior of a program by assigning abstract values to its variables, and interpreting a version of the program where all operators are considered to work on the abstract values of the variables, rather than concrete values they would assume during an execution. Such modeling can help to prove the correctness of programs.

Without going into details, let us demonstrate the usefulness of abstract interpretation on an example. Going back to the heuristic detection pattern already discussed previously (see example 2)

```
cmp  word [???+18], 10b
```

1 jnz ???

2 We already saw one way to evade heuristic detection by hiding the constant 10b. Another
3 way could be to frame the value at offset 18 from above and from below using two successive
4 comparisons.

5 cmp word [ebx+18], 10a

6 jbe dont_infect

7 cmp word [ebx+18], 10c

8 jae dont_infect

9
10 When control reaches the point after the “jae,” the word at offset 18 is both greater than
11 10a and less than 10c; therefore, it is 10b. To detect it automatically and simplify the code, we
12 can use an abstract interpretation where variables assume abstract values that are intervals of
13 numbers. If the abstract variable x has the abstract value [3..14] at one point in the program, it
14 means that the real variable x can have a concrete value only between 3 and 14 at this point of the
15 program during any execution of the program.

16
17 We are interested in the abstract value of the word at [ebx+18], so we will annotate the
18 instructions above with the abstract value of this word. On entry into the first comparison, we
19 know nothing about the word, so we will assume it can take any value, that is, its abstract value is
20 the interval [0..ffff]. The same is true on entry into the “jbe.”

21 cmp word [ebx+18], 10a ; [0..ffff]

22 jbe dont_infect ; [0..ffff]

23
24 On entry into the second comparison, the “jbe” branch has not been taken, which reduces
25 the possible range for the word to a smaller interval.

26 cmp word [ebx+18], 10c ; [10b..ffff]

27 jae dont_infect ; [10b..ffff]

28

1 ; $[10b..ffff] \cap [0..10b] = [10b..10b]$

2 Finally, on entry into the instruction following the “jae,” since the second conditional
3 jump has not been taken, the word at $[ebx+18]$ can only be in interval $[0..10b]$. Since we already
4 know it is in interval $[10b..ffff]$, the word can only have value 10b.
5

6 After determining this equality, we can introduce a piece of code that makes this assertion
7 explicit in the form of an extra conditional jump that we know can never be taken. We
8 deliberately choose the “dont_infect” label as the destination of this conditional jump, to create
9 optimization opportunities. The resulting code is:

```
10 cmp word [ebx+18], 10a  
11 jbe dont_infect  
12  
13 cmp word [ebx+18], 10c  
14 jae dont_infect  
15  
16 cmp word [ebx+18], 10b  
17 jne dont_infect
```

18 We can then apply a simplification rule to the control flow graph of the program. If two
19 back-to-back conditional branching statements have no side effects, the same destinations and one
20 of the conditions implies the other, the weaker of the two conditions may not be tested, and the
21 corresponding conditional branch instruction removed without changing the meaning of the
22 program. In this example, the condition $(\text{word}[ebx+18] \neq 10b)$ implies that $(\text{word}[ebx+18] \geq$
23 $10c)$. Therefore, we can remove the second comparison and the jump.

```
24 cmp word [ebx+18], 10a  
25 jbe dont_infect  
26  
27 cmp word [ebx+18], 10b  
28 jne dont_infect
```

1 Likewise, the first test is weaker than the second, so after applying the same rule once
2 more, we are left with the original pattern that will trigger the heuristic:

```
3  cmp  word [ebx+18], 10b  
4  jne  dont_infect
```

5 The constant folding optimization described earlier can also be seen as an abstract
6 interpretation.

7 *Program specialization*

9 Program specialization studies transformations that can be made to a program when some
10 parts of the execution context of the program are known. A special case of program
11 specialization is instruction specialization.

12 An example of instruction specialization is:

```
13       add ebx,eax → add ebx, 1234
```

14 The context of the program includes, for instance, the arguments that the program takes.

15 Consider the following program that takes three arguments:

16 Program P taking arguments i, j, k

```
17       if (i > j)  
18       print k + 2;  
19       else  
20       print i + j
```

21 The specialization of P in the context where argument i = 2 is:

22 Program P' taking arguments j, k

```
23       if (2 > j)  
24       print k + 2;  
25       else
```


1 print 2 + j

2 The specialization of P in the context where argument i = 2 and j = 1 is

3 Program P'' taking argument k

4 print k + 2;

5 At the assembly instruction level, the constant folding and copy propagation techniques
6 described earlier are in fact specialization. Thus, when we replace the following sequence of
7 instructions:

9 mov eax, 2

10 mov ebx, ecx

11 add [esi+eax], ebx

12 with the simpler sequence

13 mov eax, 2

14 mov ebx, ecx

15 add [esi+2], ecx

16 We will say that we have specialized the arguments of the "add," and that we have
17 specialized the instruction itself, based on the contextual information provided by the instructions
18 that precede it.

19 Another kind of instruction specialization is illustrated in the following example. We can
20 specialize the instruction (taken from Win32/Zmist.A)

21 xchg esp, esp

22 into a nop instruction, thus emptying its definitions set and making it a candidate for dead code
23 elimination.

24

25

26

II. Architecture of an optimizer 39

Figure 4 illustrates the overall method of optimization 40. The method begins at step 41, then an iteration loop 42-44, 49 is performed, and then the malicious code detection protocol is performed at step 45. The iteration loop comprises performing a forward pass 42, performing a backward pass 43, performing an optional code motion step 49, and performing a control flow graph reduction 44. The loop 42-44, 49 is iterated for a preselected number of iterations. Alternatively, the iteration of the loop 42-44, 49 is terminated once it is observed that there were no optimizations of the computer code performed in the most recent iteration of the loop 42-44, 49.

Figure 5 illustrates details of the forward pass procedure 42, in which at least one of the steps of Figure 5 is performed. The method begins at step 51. A peephole optimization is performed at step 52. Constant folding is performed at step 53. Copy propagation is performed at step 54. The constant folding of step 53 and/or the copy propagation of step 54 can be local and/or global. Typically, local constant folding 53 and/or copy propagation 54 is performed and, if the local techniques result in code 30 simplification, global techniques are then also performed. Forward computations related to abstract interpretation are performed at step 59. Instruction specialization is performed at step 55, and the method ends at step 56.

Figure 6 illustrates one embodiment for implementing the backward pass 43 procedure, in which at least one of the steps of Figure 6 is performed. The method begins at step 61. Backward computations related to abstract interpretation are performed at step 68. Local dead code elimination is performed at step 62. Step 63 (global dead code elimination) is optional. The decision to perform step 63 can be based upon the results of step 62, e.g., if step 62 resulted in code 30 simplification, step 63 is performed. The method ends at step 64.

Figure 3 illustrates apparatus that can execute the steps that have been discussed above. State tracking module 33 contains information concerning the status of registers, flags, different areas of memory, stacks, heaps, and state of the operating system. Peephole optimizer 31 interrogates state tracking module 33 regarding the state of the registers, flags, etc. In one embodiment, peephole optimizer 31 contains instruction reordering module 32, which receives the input instruction stream 30, creates therefrom a directed acyclic graph (such as illustrated in Figure 7), and outputs the instructions in a way that the instructions that are likely to be peephole optimized 52 by remaining portions of the peephole optimizer 31 are next to each other.

Virtual state memory module 35 gives the state of the registers, flags, etc., at each stage of the instruction stream 30. State tracking module 33 is the interface between virtual state memory module 35 and peephole optimizer 31, instruction specialization module 34, and driver module 36.

State tracking module 33 provides input for all of the major steps of the optimization 40.

Driver module 36 performs all of the optimization 40 steps except for peephole optimization 52 and program specialization 55.

Symbolic instruction module 38 holds symbolic representations of processor instructions, typically a set of nodes in the shape of a control flow graph.

The user can provide inputs to the optimization 40 by means of providing initial conditions to state tracking module 33. That gives one the ability to optimize when it would not otherwise be possible, e.g., in cases where the instruction stream 30 contains a buggy virus. For example, the user may conclude by observing the behavior of the virus that certain instructions referencing a certain memory range are dead; and the user then provides this information to state tracking module 33.

Considerations on code transformations

1 During the presentation of the optimization techniques 40 above, we voluntarily skipped
2 over some conditions that are verified in order for the code transformations to be correct. We now
3 revisit some problematic aspects of these techniques in finer detail.

4 Consider the peephole optimization 52 that transforms a pair of back-to-back push and
5 pop instructions into a mov instruction. The original code may look like the following (taken
6 from Win32/Simile.A)

```
7 push      dword [40023fb0]  
8  
9 pop       eax
```

10 It seems safe to simplify this pair of instructions into one mov:

```
11 mov      eax, dword [40023fb0]  
12
```

13 While this transformation (a typical peephole optimization 52) would usually be correct,
14 there are also some special contexts where it is not, among which:

- 15 1. If the stack value below the stack pointer is used after the pop.
- 16 2. If the access to the memory location [40023fb0] causes an exception.
- 17 3. If the stack pointer used by the push instruction is pointing to the pop instruction (that is,
18 the instruction sequence is self-modifying).
- 19 4. If the processor is in tracing mode and an interrupt occurs after every instruction.

20 All of these special contexts could be used as anti-debugging tricks. Win32/Chiton.E
21 (a.k.a. Win32/Efish) checks the value below the stack pointer to see if it has been modified due to
22 a debugger. Some viruses use the Structured Exception Handling mechanism of Windows to
23 transfer control and thus make emulation and analysis more difficult (Win32/Magistr,
24 Win32/Efortune, Win32/Hezhi, Win32/Chiton). Self-modifying code is very common in viruses
25 (all polymorphic viruses 10 decrypt their own code 12). Win32/Perenast executes applications in
26 tracing mode to implement Entry-Point Obscuring. The decompression code of the tELock
27
28

1 executable packer runs in tracing mode and keeps count of the number of instructions executed,
2 and then verifies it is below a threshold to ensure no debugger is present.

3 Drawing from these observations, we should make sure that the context of the push/pop
4 pair is proper before optimizing the pair.

- 5 1. Live variable analysis should tell if the stack value below the stack pointer is dead on exit
6 from the pop instruction. This is very often easy to prove if the stack is reused later in the
7 code, since any push will kill this value.
- 8 2. Instruction specialization 55 according to constant folding 53 and copy propagation 54
9 should indicate if the argument of the push is likely to trigger an exception.
- 10 3. Constant folding 53 and copy propagation 54 should indicate if the stack pointer was
11 earlier set to point to the code.
- 12 4. Analysis of earlier code should reveal if the trap flag of the processor has been set and the
13 processor is in tracing mode when the push/pop sequence is reached.

14 Of course, the four problems stated above are impossible to solve perfectly (theoretically
15 they are all undecidable). In practice, however, there is a good chance that if the code preceding
16 the push/pop pair explicitly attempts to set up a wrong memory location as the push argument, or
17 to point the stack pointer to the instructions, a code analysis using constant folding 53 and copy
18 propagation 54 would reveal this fact. In the absence of a flagrant sign of such manipulations, the
19 optimization 40 can be done assuming the simplest context.

20 When optimizing polymorphic virus code 10, best effort is often enough. Optimizing
21 towards exactly equivalent code is a desirable property, for instance to ensure that the emulation
22 of optimized code 37 will yield proper results, but not a necessity as long as the output 37 of the
23 optimizer 39 can be used reliably for pattern matching, checksumming, heuristics, and other kinds
24 of information gathering related to virus detection.

1 The push/pop example suggests that it is preferable to do at least some part of the
2 peephole optimization 52 after the constant folding 53 and copy propagation 54. However, we
3 said earlier that local constant folding 53 was improved if peephole optimization 52 was used for
4 fake conditional jumps removal. To overcome this dilemma, in one embodiment, there are two
5 peephole optimizer steps 52, one that runs as the first step during the decoding of the machine
6 instructions 30, and one that operates later, when some data-flow analysis 53,54 has already been
7 done. In fact, we can use the same peephole optimizer 31 in several iterations of the loop 42-44,
8 49.

10 Another example that illustrates the usefulness of doing live variable analysis before
11 peephole optimization 52 is the application of algebraic identities on back-to-back logic or
12 arithmetic instructions. When consecutive instructions have the same destination argument and a
13 constant source argument, some simplifications may be possible.

15 The following two instructions (from Win32/Simile.A)

16 and ebx, bfadfffe

17 and ebx, 6efbffffd

18 can always be optimized to:

19 and ebx, 2ea9fffc

20 where the new mask on the right-hand side is the bitwise “and” of the two original masks. The
21 optimization 40 is possible regardless of the context because the flags produced by the second
22 “and” of the instruction pair are the same as the flags produced by the optimized “and” in all
23 cases.

25 On the contrary, the following two instructions:

26 add ebx, 2

27 add ebx, 2

cannot, without some context information, be optimized safely to:

```
add ebx, 4
```

because the resulting carry flag may differ (consider a case where `ebx = ffffffff` on entry into the instruction pair.) If previous live variable analysis revealed that the flags are dead after the second “add,” the optimization 40 is proper.

Less obvious algebraic identities cannot be detected by a peephole optimizer 31, because they require reordering the terms of expressions. Consider the following example:

```
mov ecx, eax
```

```
and eax, ebx
```

```
not ebx
```

```
and ecx, ebx
```

```
or ecx, eax
```

Whatever the value of register `ebx`, `ecx` on exit is a copy of `eax` on entry.

Dependency DAG construction and reordering of instructions

One limitation of a simple peephole optimizer 31 is that it does not naturally handle optimizations of non-contiguous instruction sequences. Consider the following example:

```
(I1) push eax
```

```
(I2) and ebx, ff
```

```
(I3) pop ecx
```

```
(I4) and ebx, ff00
```

```
(I5) add ebx, ecx
```

Furthermore, let us assume that the flags and stack are dead on exit from the final “add.” Under these conditions, it should be obvious that a first optimization step for this block of code

would be to change the push/pop pair into a mov instruction, and to combine the two “and” instructions together:

```
mov    ecx, eax
and    ebx, 0
add    ebx, ecx
```

From there, copy propagation 54, instruction specialization 55, and dead code elimination 62 easily lead to:

```
mov    ecx, eax
mov    ebx, eax
```

Unfortunately, the first optimization step is out of reach for a simple peephole optimizer 31, because none of the pairs of contiguous instructions in the original block can be combined. The problem resides in the intertwined sequences of instructions belonging to parallel dependency chains. To solve this problem, peephole optimization 52 can be applied to the output of a filter 32 that reorders the instructions.

When processing a block of instructions, we build a directed acyclic graph (DAG) where the nodes represent instructions and the edges represent a dependency relationship between the instructions. More exactly, an edge from A to B indicates that some definitions of instruction B reach instruction A and are either used or killed by instruction A. The DAG of the original block above is illustrated in Figure 7.

Paths of the DAG express the dependency chains between instructions. For instance, instruction 5 must come after both instruction 3 and instruction 4, because it uses results produced by both these instructions. Instruction 3 must come after instruction 1, and instruction 4 must come after instruction 2.

1 Having built this DAG structure describing all instructions of a block, we can create an
2 equivalent block by visiting the nodes of the DAG and emitting their instructions in postorder,
3 that is, emitting a node by instruction reordering module 32 within peephole optimizer 31, only
4 after all the nodes it points to have been emitted already. The most recently emitted instruction is
5 the first instruction in the block under construction, i.e., the block is created bottom to top.
6

7 There are multiple solutions to this problem because, at any moment during the emission
8 of the instructions, there might be multiple available nodes whose descendants have all been
9 emitted. In such a case, we break ties by picking an available node that offers a peephole
10 optimization 52 opportunity with the most recently emitted instruction, if such a node exists.
11 Following the algorithm, the resulting block for the example above exposes the peephole
12 optimization spots quite nicely:

13 (I1) push eax

14 (I3) pop ecx

15 (I2) and ebx, ff

16 (I4) and ebx, ff00

17 (I5) add ebx, ecx

18
19 The algorithm can be extended to handle cases when a peephole optimization 52 would
20 lead to the creation of new opportunities, like the case of nested push/pop pairs. The choice of
21 available nodes during code emission can also be dependent on other criteria than just peephole
22 optimization. Picking the emitted instructions based on an ordering of the opcodes can help
23 simplify later pattern matching in the resulting block.
24

25 *Approximation of the control flow graph*

26 The control flow of a program may depend on the data in non-trivial ways. For instance,
27 the program may contain jump tables that implement high-level switch statements. In such a case,
28

1 code addresses are part of the program data, and a data-flow analysis is required to avoid missing
2 some paths in the control flow.

3 Jump tables occur naturally in compiled high-level language programs, but some other
4 issues are (almost always) specific to programs written in assembly language, like self-modifying
5 code or idiomatic use of some instruction sequences. One example is the call/pop sequence that
6 appears very frequently in viruses. It can be used to obtain a pointer to some data embedded in the
7 code, in which case the call should really be handled as a jump, because it never returns. Another
8 example is the push/ret sequence that can be used to jump to an absolute address.

9
10 Given a program written in a high-level language, it is easy to overestimate its possible
11 control flow paths, whereas it is hard to do so for a virus because of call/pop and push/ret
12 sequences whose control flow approximation already requires some data-flow analysis.

13
14 An iterative approach may be appropriate, where control flow is first estimated
15 heuristically by tracing the code and applying some reasonable rules (calls always return,
16 exceptions do not occur), and then some data-flow analysis and optimization takes place. Then,
17 based on the results of the data-flow analysis (steps 42, 43, 49), some control flow paths are
18 added and some are removed (step 44). Finally, parts of data-flow analysis and optimization
19 results are invalidated, and recomputed in the next pass of the iteration loop 40.

20 *Reduction of the control flow graph*

21
22 Once dead code elimination 62 has removed useless instructions from basic blocks and
23 code motion 49 has moved instructions across block boundaries, some blocks may turn out
24 empty, or almost empty.

25 If a block is empty, except maybe for a last unconditional branch, the control flow can be
26 modified 44 so that predecessors of the block branch directly to the successor of the block, and
27 the block can be removed.

1 If a block ends with a conditional branch to itself (the block is a loop), and if all
2 instructions left in the block only determine the outcome of the branch, the block is a dummy
3 loop and may be removed 44. Here is an example of a dummy loop taken from a sample of
4 Win95/Zexam:

5
6 101704a:

```
7     shrd eax, edx, 17
8     imul ecx, ecx, ecx
9     inc  eax
10    sub  esi, a81a9913
11    mov  eax, ecx
12    imul ebx, ebx, ebx
13    add  ebp, b3c0136a
14    bsr  ebx, ecx
15    btr  ebx, 1f
16    not  ebx
17    mov  ecx, 11ece82
18    cmp  esi, f5b744be
19    jnz  101704a
```

21 On exit from the loop, the processor flags and registers eax, ebx, ecx and ebp are dead
22 (they are killed by the code following the loop). Global dead code elimination 63 yields the
23 following code:

24
25 looptop:

```
26     sub  esi, a81a9913
27     cmp  esi, f5b744be
```

```
1      jnz  looptop
2      mov  esi, f5b744be
```

3 The control flow graph for this code is illustrated in Figure 9(a). The assignment to
4 register esi inserted after the loop does not change the meaning of the program, since it is
5 redundant with the exit condition of the loop. This optimized loop now contains only instructions
6 that affect its conditional branch, since the flags and esi are dead on exit. Therefore, the loop can
7 be removed. (We assume that the loop exits at some point; in other words, it is not an infinite
8 loop. Some heuristics can help in determining this.) The control flow graph for this code after
9 loop removal is illustrated in Figure 9(b).

11 As a result of dummy loops elimination, emulation of polymorphic decryptors 11 can
12 become much faster, especially if loops can be nested.

14 Another useful reduction 44 of the control flow graph is the elimination of calls to blocks
15 that contain a single “ret” instruction.

16 *Specifying boundary conditions*

17 Two types of information participate in the resolution of data-flow equations: data
18 gathered from the nodes of the control flow graph (the basic blocks), and boundary conditions
19 that apply on the start and exit nodes of the control flow graph. For instance, live variable analysis
20 is a backwards analysis that propagates information up through the basic blocks. For the last basic
21 block of a program (in execution order), it is customary to assume that all variables are dead on
22 exit from the block. This boundary condition expresses the fact that no variables are ever going to
23 be used after the program exits.

25 Boundary conditions are not so clear-cut in the case of programs containing self-
26 modifying code. In a polymorphic virus 10, the decryptor 11 produces a piece of code 12 and then
27
28

1 executes it 12. The set of live variables on exit from the decryptor 11 is hard to determine,
2 because it depends on the register and memory usage of the code 12 it decrypts.

3 To be conservative, one can assume that all variables are live on exit from the decryptor
4 11, but it could lead to inefficient optimization in some cases. Another possibility is to guess that
5 some variables are dead, optimize the decryptor 11 based on this assumption, emulate the
6 resulting code 12, and then verify that the variables are actually dead by analyzing the decrypted
7 code.
8

9 Rather than guessing boundary conditions, an alternative is to let a user specify them to
10 the state tracking module 33 of the optimizer 39. More generally, allowing the user to specify
11 conditions at various program points makes the optimizer 39 more flexible, and capable of
12 handling buggy code produced by some polymorphic engines 10. Win32/Hezhi sometimes fails to
13 finish its decryption loop 11 with a proper backwards jump. Win32/Simile.D produces some
14 corruptions where the polymorphic decryptor 11 patches itself. User-supplied options would
15 allow the optimizer 39 to circumvent these problems.
16

17 Compared with tracing, emulation, and X-raying, code optimization 40 can do one thing
18 that none of these other techniques can, namely simplify code 30. Being able to work on readable
19 code when analyzing the body 22 of a metamorphic virus 20 can be a tremendous help (see, e.g.,
20 Illustration D on Win95/Puron). Optimization 40 also makes exact identification of metamorphic
21 virus 20 variants possible, based on their simplified body 22. Variant identification is an
22 advantage for multiple reasons.
23

24 We use the term "tracing" to refer to the technique that consists in doing a partial
25 disassembly of a program and attempting to follow its control flow based on simple rules.
26 Typically, in tracing, only the length of instructions is calculated, except for branches that must
27 be fully disassembled to follow them.
28

1 Tracing can be used to detect polymorphic decryptors 11 that present some easily
2 recognizable characteristics, but are split into islands of code linked together by branches
3 (Win32/Marburg, Win32/Orez.) It can also been used to detect metamorphic bodies 22 that use a
4 weak form of metamorphism where some fixed instructions are always present.

5
6 The first phase of an optimizer 39 is instruction decoding, which is very similar to tracing
7 in spirit. An optimizer 39 is slower than a tracer because of the extra work associated with full
8 instruction decoding. However, it is usable in more situations, for instance when the code 30
9 contains indirect jumps through registers whose values are built dynamically. An efficient hybrid
10 approach would be to simply trace the code 30 and check some decryptor characteristics up to a
11 point where such a problematic indirect branch is used; then do a complete instruction decoding,
12 followed by a data-flow analysis 42, 43, 49 on the subset of instructions that contribute to the
13 branch destination (this subset is called a program slice).

14
15 Previous paragraphs already discussed several ways to make emulation faster by
16 optimizing 40 the code 30 to emulate. In many situations, pattern matching on the optimized code
17 can also replace emulation for the purpose of detection (see the below Illustrations), though
18 emulation may still be needed for exact variant identification. For very complex polymorphic
19 viruses 20, the emulation speed can be improved by factors of hundreds.

20
21 Systematically optimizing 40 code before emulating it results in a performance hit if the
22 original code 30 is already as simple as it can be. However, the slowdown is by a small constant
23 ratio. If local optimizations are used first and global optimizations take place only if local
24 optimizations gave some improvements, the extra time is linear in the code 30 size. This is
25 unlikely to be a problem, compared for instance to the cost of input/output.

26
27 As to X-raying, which is a technique that performs a known cleartext attack on the
28 encrypted virus body, it might be replaced by optimization 40 when X-raying is used, because

1 emulation of the decryptor would take too long, or when emulation is not an option because the
2 virus produces buggy decryptors. Emulation of the optimized decryptor, or pattern matching on it,
3 may be a viable alternative.

4 If X-raying is used because the virus uses Entry-Point Obscuring and the location of the
5 decryptor is unknown (or, at least, not easily guessable), optimization 40 may not be able to help.

6 *Dead code elimination as a heuristic*

7
8 Another use of optimization 40 is as a heuristic to detect polymorphic code 10. Most
9 polymorphic engines 10 produce many redundant instructions, whereas a typical program has
10 almost no dead code.

11 There are a few exceptions where dead code can be useful in a normal program. The use
12 of nop instructions to allow pairing of instructions on superscalar processors, or to align loop top
13 addresses on even boundaries can speed up execution. Dummy memory reads whose results are
14 discarded are sometimes used to prefill the processor cache. Likewise, some processor
15 instructions, like "pause" and other processor hints, are functionally dead but affect how the
16 program runs.

17
18 However, the amount of dead code in the cases described above represents a very small
19 percentage of the overall program. On the other hand, the dead code ratio in the output of
20 polymorphic engines 10 is typically higher than 25%, and sometimes much more (see some
21 examples in the below Illustrations.)

22
23 The presence of dead code by itself is not enough to declare a program viral, since
24 polymorphic code 10 exists in legitimate executables, such as packed files (Aspack), but it is
25 suspicious enough to warrant further investigation. Therefore, a method embodiment of the
26 present invention comprises performing a dead code elimination procedure on the computer code
27 30; noting the amount of dead code eliminated during the dead code elimination procedure; and
28

1 when the amount of dead code eliminated during the dead code elimination procedure exceeds a
2 preselected dead code threshold, declaring a suspicion of malicious code in the computer code 30.

3 **III. Illustrations**

4 The data presented here were obtained by running a prototype optimizer 39 containing
5 most of the modules described above on some code samples of polymorphic 10 and metamorphic
6 20 viruses. In each case, we list the disassembly of the original code 30, followed by the output of
7 the optimizer 39.
8

9 **Illustration A**

10 *Win95/Zperm.B*

11 Win95/Zperm is a metamorphic virus 20 that permutes its body 22. This example shows part of
12 the API resolution routine, before and after jump removal.
13

14 Original code:

```
15 4118db: stosd
16 4118dc: mov     eax, ae17c571
17 4118e1: call    edx
18 4118e3: jmp     41b65b
19 418184: mov     eax, 1fc0eae     ; entry-point
20 418189: call    edx
21 41818b: jmp     4118db
22
23 418534: stosd
24 419657: mov     eax, 7b4842c1
25 41965c: call    edx
26 41965e: stosd
27 41965f: mov     eax, 32432444
28
```



```
1 419664: call     edx
2 419666: jmp      418534
3 41b65b: stosd
4 41b65c: jmp      419657
```

Optimized code:

```
6
7     mov     eax, 1fc0eae
8     call    edx
9     stosd
10    mov     eax, ae17c571
11    call    edx
12    stosd
13    mov     eax, 7b4842c1
14    call    edx
15    stosd
16    mov     eax, 32432444
17    call    edx
18    stosd
```

Since the calls are in order in the optimized code, a simple search string can be used to detect the virus 20.

Illustration B

Win95/Zmorph

Win95/Zmorph is a polymorphic virus 10 that builds its body 12 on the stack. This example illustrates constant folding 53.

Original code:

1	4122a7:	mov	ebx, d1632349
2	4122ac:	mov	edx, 38d9cdd5
3	4122b1:	add	ebx, 810ad92a
4	4122b7:	mov	esi, dcf4a826
5	4122bc:	rol	edx, b
6	4122bf:	sub	esi, 4c641727
7	4122c5:	xor	edx, 8963fd03
8	4122cb:	add	ebx, ad8ddd76
9	4122d1:	mov	eax, 38c30f5d
10	4122d6:	mov	ecx, dded6aa9
11	4122db:	not	ecx
12	4122dd:	sub	eax, 77b356f7
13	4122e2:	mov	edi, 4c618901
14	4122e7:	bts	edi, b
15	4122eb:	add	edi, 8833c388
16	4122f1:	ror	ecx, 15
17	4122f4:	push	esi
18	4122f5:	push	ebx
19	4122f6:	bswap	edx
20	4122f8:	push	eax
21	4122f9:	xor	esi, ecx
22	4122fb:	xor	eax, 1592fcef
23	412300:	imul	ebx, ebx, 30e081f5
24	412306:	push	edi

```
1 412307: bts      esi, b
2 41230b: add      edi, f42bc34b
```

3 Optimized code:

```
4      push      909090ff
5      push      fffbd9e9
6      push      c10fb866
7      push      d4954c89
8
9      mov       edi, d4954c89
10     add       edi, f42bc34b
11     mov       eax, d49d4489
12     mov       ecx, 94aab110
13     mov       edx, c5540d47
14     mov       ebx, 40b5f4fd
15     mov       esi, 43a29ef
16     mov       edi, c8c10fd4
17
```

18 The four highlighted pushes create the tail of the virus 10, and they can be used for
19 detection. The movs and the add reflect the processor state at the end of block.

20 Illustration C

21 *Win95/Zmist.A*

22
23 Win95/Zmist is a metamorphic and entry-point obscuring virus 20. This example illustrates
24 constant folding 53. (The entry-point of the virus body 22 was given as a parameter to the
25 optimizer 39.)

26 Original code:

```
27 404945: jmp      40494a
28
```

1	40494a:	pusha	
2	40494b:	xor	eax, eax
3	40494d:	sub	eax, 87868600
4	404952:	push	eax
5	404953:	xor	eax, 7274542e
6	404958:	push	eax
7	404959:	add	eax, 245f3e33
8	40495e:	push	eax
9	40495f:	xor	eax, 48181f08
10	404964:	push	eax
11	404965:	sub	eax, 19540004
12	40496a:	push	eax
13	40496b:	mov	esi, esi
14	40496d:	xor	eax, 204f1045
15	404972:	push	eax
16	404973:	mov	eax, eax
17	404975:	add	eax, f9ff064e
18	40497a:	push	eax
19	40497b:	xor	eax, 1501044e
20	404980:	push	eax
21	404981:	sub	eax, 9fb03a9
22	404986:	push	eax
23	404987:	push	esp
24	404988:	push	d0498cd4

28

1 Optimized code:

2 pusha

3 push 78797a00

4 push a0d2e2e

5 push 2e6c6c61

6 push 66747369

7 push 4d207365

8 push 6d6f6320

9 push 676e696e

10 push 726f6d20

11 mov eax, 726f6d20

12 sub eax, 9fb03a9

13 push 68746977

14 push esp

15 push d0498cd4

16 mov eax, 68746977

17 The data pushed on the stack is a text that reads “with morning comes Mistfall...” and can
18 be used for detection. The movs and add that are left would be removed by global dead code
19 elimination 63 if the analysis context was extended to include the code following this snippet.

20 Illustration D

21 Win95/Puron

22 Win95/Puron is a metamorphic virus 20 that mixes dead code with the meaningful
23 instructions of its body 22, and splits its body 22 into islands of code linked by jumps.

1 This example is taken from the routine that searches the address base of the kernel module
2 in memory. It illustrates dead code elimination and jump removal.

3 Original code:

```
4        40a3a5: lea        esi, [edi+62309cc]
5
6        40a3ab: pop        ebx
7
8        40a3ac: jnz        40aa2f
9
10       40a3b2: lea        esi, [edi+3627dfc]
11
12       40a3b8: push       ecx
13
14       40a3b9: sub        ecx, 400
15
16       40a3bf: pop        ecx
17
18       40a3c0: mov        ebp, 6626b32
19
20       40a3c5: jmp        40a517
21
22       40a517: mov        bh, dh
23
24       40a519: movsx      ebp, bh
25
26       40a51c: jmp        40aa1a
27
28       40a5d6: dec        edx               ; entry-point
29
30       40a5d7: mov        ebp, 2ee8d12
31
32       40a5dc: jmp        40abf9
33
34       40a6e8: mov        ecx, dword [edx+3c]
35
36       40a6eb: mov        ebx, ebp
37
38       40a6ed: mov        esi, 4f5ce1f
39
40       40a6f2: mov        bh, b1
41
42       40a6f4: cmp        word [edx], 5a4d
43
44       40a6f9: push       ecx
```

1	40a6fa:	mov	ebx, ebp
2	40a6fc:	lea	esi, [edi+3fee834]
3	40a702:	jmp	40a3a5
4	40a756:	pop	eax
5	40a757:	mov	ebx, ebp
6	40a759:	mov	esi, 4b5d687
7			
8	40aa1a:	jbe	40aa28
9	40aa1c:	xor	ecx, ecx
10	40aa1e:	mov	bh, e1
11	40aa20:	lea	ebp, [edx+7c50c63]
12	40aa26:	mov	edi, esi
13	40aa28:	cmp	dword [edx+ecx], 4550
14			
15	40aa2f:	popa	
16	40aa30:	mov	ebx, edx
17	40aa32:	mov	esi, 70b62af
18	40aa37:	jnz	40a5d6
19	40aa3d:	jmp	40aad6
20	40aab4:	pusha	
21	40aab5:	jmp	40a6e8
22			
23	40aad6:	pop	dword [0]
24	40aae3:	mov	ebx, ebp
25	40aae5:	mov	bh, dh
26	40aae7:	mov	ebx, 5b2b5d8
27	40aaec:	lea	edi, [ebp+65e63a2]
28			

```

1      40aaf2:  jmp          40a756
2      40abf9:  xor          edi, 78f710c
3      40abff:  mov          ebx, 64891f8
4      40ac04:  mov          al, bh
5      40ac06:  mov          ecx, eax
6
7      40ac08:  jmp          40aab4
8  Optimized code:
9  block 0
10     dec      edx
11     xor      edi, 78f710c
12     mov      al, 91
13     mov      ecx, eax
14     mov      ebx, 64891f8
15     mov      ebp, 2ee8d12
16
17     pusha
18     mov      ecx, dword [edx+3c]
19     cmp      word [edx], 5a4d
20     push     dword [edx+3c]
21     pop      ebx
22
23     jnz      2          ; destinations are block numbers
24  block 1
25     push     ecx
26     sub      ecx, 400
27     pop      ecx
28

```



```

1      jbe      5
2  block 4
3      mov      ecx, 0
4  block 5
5      cmp      dword [edx+ecx], 4550
6  block 2
7      popa
8      mov      esi, 70b62af
9      jnz      0
11 block 3
12      pop      dword [0] ; an fs: selector is missing
13      lea      edi, [ebp+65e63a2]
14      pop      eax
15      mov      ebx, ebp
16      mov      esi, 4b5d687

```

18 The highlighted instructions are dead code that remains because of the pusha instruction in
 19 block 0. Pusha uses all registers, which is why the register assignments preceding it seem
 20 necessary. In fact, the pushed registers are later popped in block 2 and discarded. This “tunnel
 21 effect” can be avoided by using a fine-grained live variable analysis on the stack elements.

22 Notice also the presence of a push/pop sequence in block 0. The sequence was not
 23 peephole-optimized 52 into a mov, because the two instructions are separated by dead
 24 instructions in the original code, and the peephole optimization 52 took place before dead code
 25 elimination 62. As a result, even though ebx is dead after the “pop ebx” because it is killed by the
 26 popa instruction later, the push/pop pair remains because of its use of the stack. The prototype
 27
 28

optimizer 39 used in this Illustration does not implement the dependency DAG construction described earlier, which would resolve this problem.

Illustration E

Win32/Dislex

Win32/Dislex is a complex polymorphic virus 10 based on the Lexotan engine.

This example is taken from the polymorphic loop 11 that decrypts the data area 12 of the virus 10.

Once decrypted, the content of the data area 12 can be used for detection. This example illustrates the use of optimization 40 to speed up emulation.

Original code:

```
4030ca: pusha
4030cb: jmp      4041c2
4032ef: add     edx, ebx
4032f1: inc     edi
4032f2: movzx   edi, dl
4032f5: jmp     403809
403728: jnz     406d35
40372e: mov     edi, 7ce07ac
403733: mov     edi, ebp
403735: movzx   edi, dl
403738: jmp     408841
4037cb: push    eax      ; entry-point
4037cc: jmp     4030ca
403809: mov     dword [esi+ffffffc], eax
40380c: lea     edi, [ebp+7f9a292]
```

1	403812:	jmp	406ff5
2	403e90:	mov	edx, dword [40947e]
3	403e96:	mov	edi, ebp
4	403e98:	jmp	406ef7
5	4041c2:	lea	ebp, [edx+5a5f84b]
6	4041c8:	mov	eax, ecx
7	4041ca:	mov	di, ab04
8	4041ce:	mov	ah, dh
9	4041d0:	movzx	ebp, al
10	4041d3:	movsx	edi, dx
11	4041d6:	or	edi, 76d9ecc
12	4041dc:	lea	eax, [ecx+5e4f6]
13	4041e2:	and	ah, ce
14	4041e5:	jmp	404780
15	404780:	mov	esi, 4091ca
16	404785:	lea	eax, [ecx+64f77a6]
17	40478b:	mov	ah, 32
18	40478d:	add	ah, 8a
19	404790:	add	ah, e2
20	404793:	mov	ah, 2e
21	404795:	mov	ah, dh
22	404797:	sub	eax, 5731a19
23	40479d:	push	ad
24	4047a2:	lea	ebp, [edx+56dfddb]

1	4047a8:	lea	edi, [ebp+2785942]
2	4047ae:	mov	eax, 4e1bb89
3	4047b3:	lea	ebp, [edx+52613cb]
4	4047b9:	lea	edi, [ebp+2dd96f2]
5	4047bf:	mov	eax, 4b398f9
6			
7	4047c4:	inc	edi
8	4047c5:	mov	ah, dh
9	4047c7:	mov	ah, dl
10	4047c9:	or	edi, 707681c
11	4047cf:	adc	ah, c6
12	4047d2:	jmp	405e2b
13	405e2b:	pop	ecx
14	405e2c:	sbb	eax, 25d07d9
15	405e32:	mov	edi, ebp
16	405e34:	mov	eax, 246d911
17	405e39:	sub	eax, 2029949
18	405e3f:	cmp	ebp, 54ea55a
19	405e45:	movsx	eax, bh
20	405e48:	mov	bp, 85b2
21	405e4c:	jmp	403e90
22	406d35:	lodsd	
23	406d36:	or	edi, 7bb6e04
24	406d3c:	mov	edi, ebp
25	406d3e:	sbb	edi, 7586034
26			
27			
28			

1	406d44:	movzx	edi, dx
2	406d47:	lea	edi, [ebp+63d582]
3	406d4d:	lea	edi, [ebp+3292da]
4	406d53:	xor	eax, edx
5	406d55:	mov	di, 894
6	406d59:	movzx	edi, dx
7	406d5c:	jmp	4032ef
8	406ef7:	mov	ebx, dword [409482]
9	406efd:	mov	edi, ebp
10	406eff:	mov	ax, 7029
11	406f03:	lea	edi, [ebp+2f28d72]
12	406f09:	lea	edi, [ebp+3d8c512]
13	406f0f:	or	edi, 467e90c
14	406f15:	movsx	edi, dl
15	406f18:	lea	edi, [ebp+4563c1a]
16	406f1e:	mov	edi, 4c467d4
17	406f23:	jmp	406d35
18	406ff5:	lea	edi, [ebp+10258ca]
19	406ffb:	movsx	edi, dl
20	406ffe:	mov	edi, ebp
21	407000:	mov	edi, ebp
22	407002:	movzx	edi, dx
23	407005:	mov	di, cf84
24	407009:	mov	edi, ebp

28

```

1    40700b:  mov     di, 21b4
2    40700f:  mov     di, f34c
3    407013:  jmp     407d1b
4    407d1b:  dec     ecx
5    407d1c:  lea     edi, [ebp+7709302]
6    407d22:  movzx   edi, dl
7    407d25:  jmp     403728
8    408841:  lea     ebx, [eax+18346b1]
9    408847:  mov     ebp, edx
10

```

11 Optimized code:

12 block 0

```

13     push     eax
14     pusha
15     mov     esi, 4091ca
16     push     ad
17     pop      ecx
18     mov     edx, dword [40947e]
19     mov     ebx, dword [409482]
20

```

21 block 1

```

22     lodsd
23     xor     eax, edx
24     add     edx, ebx
25     mov     dword [esi+fffffffc], eax
26     dec     ecx
27

```

28

1 jnz 1 ; destinations are block numbers

2 block 2

3 movzx edi, dl

4 lea ebx, [eax+18346b1]

5 mov ebp, edx

6

7 The original loop 11 contains more than thirty instructions, whereas the optimized loop

8 contains six instructions. Emulating the optimized code 37 will thus speed up emulation by a

9 factor of five. In some cases, Win32/Dislex will produce loops with hundreds of dead

10 instructions, making the benefit of optimizing before emulation even greater.

11 **Illustration F**

12 *Win32/Simile.A*

13

14 Win32/Simile is a polymorphically-encrypted metamorphic virus 20.

15 This example is taken from part of a decryptor 21 that resolves the address of the VirtualAlloc

16 API function dynamically. This example illustrates copy propagation 54, constant folding 53, and

17 dead code elimination 62.

18 Original code:

19 4000b0dd: mov dword [40023380], eax

20 4000b0e3: mov edx, 416c6175

21 4000b0e8: mov ecx, edx

22 4000b0ea: push 74726956

23 4000b0ef: pop dword [4002421b]

24 4000b0f5: mov edi, dword [4002421b]

25 4000b0fb: mov dword [40023480], 99ff02a7

26 4000b0fb: mov dword [40023480], 99ff02a7

27 4000b105: xor dword [40023480], 2649b0b1

28

1	4000b10f:	xor	dword [40023480], dcd9de7a
2	4000b119:	push	dword [40023480]
3	4000b11f:	pop	dword [40023b5b]
4	4000b125:	mov	esi, dword [40023b5b]
5	4000b12b:	clc	
6	4000b12c:	lea	ebp, [esi]
7	4000b12e:	lea	ebx, [ecx]
8	4000b131:	mov	dword [40023374], ebx
9	4000b137:	mov	dword [40023370], edi
10	4000b13d:	mov	dword [40023378], ebp
11	4000b143:	lea	edi, [8abaf6b]
12	4000b149:	add	edi, 7545e095
13	4000b14f:	lea	ecx, [edi]
14	4000b151:	mov	dword [4002337c], ecx
15	4000b157:	lea	ecx, [e49e73bc]
16	4000b15d:	add	ecx, 5b63bfb4
17	4000b163:	mov	dword [400238a0], ecx
18	4000b169:	push	dword [400238a0]
19	4000b16f:	mov	eax, dword [40023380]
20	4000b175:	clc	
21	4000b176:	mov	ecx, eax
22	4000b178:	mov	dword [40024113], ecx
23	4000b17e:	push	dword [40024113]
24	4000b184:	mov	edi, 400253a8

28


```

1  4000b18a:  call     dword [edi]
2  Optimized code:
3      mov     dword [40023380], eax
4      push    40023370
5      mov     dword [40024113], eax
6      push    eax
7      mov     ebx, 416c6175
8      mov     ebp, 636f6c6c
9      mov     esi, 636f6c6c
10     mov     edi, 400253a8
11     mov     byte [40023370], 56 ; V
12     mov     byte [40023371], 69 ; i
13     mov     byte [40023372], 72 ; r
14     mov     byte [40023373], 74 ; t
15     mov     byte [40023374], 75 ; u
16     mov     byte [40023375], 61 ; a
17     mov     byte [40023376], 6c ; l
18     mov     byte [40023377], 41 ; A
19     mov     byte [40023378], 6c ; l
20     mov     byte [40023379], 6c ; l
21     mov     byte [4002337a], 6f ; o
22     mov     byte [4002337b], 63 ; c
23     mov     byte [4002337c], 0
24     mov     byte [4002337d], 0
25
26
27
28

```

```

1      mov      byte [4002337e], 0
2      mov      byte [4002337f], 0
3      mov      byte [40023480], 6c
4      mov      byte [40023481], 6c
5      mov      byte [40023482], 6f
6      mov      byte [40023483], 63
7      mov      byte [400238a0], 70
8      mov      byte [400238a1], 33
9      mov      byte [400238a2], 2
10     mov      byte [400238a3], 40
11     mov      byte [40023b5b], 6c
12     mov      byte [40023b5c], 6c
13     mov      byte [40023b5d], 6f
14     mov      byte [40023b5e], 63
15     mov      byte [4002421b], 56
16     mov      byte [4002421c], 69
17     mov      byte [4002421d], 72
18     mov      byte [4002421e], 74
19
20     call      dword [400253a8]      ; va of GetProcAddress
21
22

```

The highlighted parts can be used for pattern matching.

The optimized code 37 is longer than the original, but this is simply a consequence of expressing the memory state on exit from the block as a series of byte assignments. The flags and registers `eax`, `ecx`, and `edx` are considered dead on entry into `GetProcAddress`, which allows some dead code elimination 62. The other registers and all memory locations are considered live, to be

1 conservative, but global dead code elimination 63 across API calls could help simplify the code
2 further.

3 The above description is included to illustrate the operation of the preferred embodiments
4 and is not meant to limit the scope of the invention. The scope of the invention is to be limited
5 only by the following claims. From the above discussion, many variations will be apparent to one
6 skilled in the art that would yet be encompassed by the spirit and scope of the present invention.
7

8 What is claimed is:
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28